

An Automatic Finite-Sample Robustness Metric: Can Dropping a Little Data Make a Big Difference?

Ryan Giordano (rgiordan@mit.edu)¹
January 2022

¹With coauthors Rachael Meager (LSE) and Tamara Broderick (MIT)

Dropping data: Motivation

More data & cheaper computation \Rightarrow

Statistical analyses are playing larger roles in decision making.

Decisions are important: We want **trustworthy** conclusions.

Data / models not always perfect: We want **robust** conclusions.

Would you be concerned if you could **reverse your conclusion** by removing a **small proportion** (say, 0.1%) of your data?

Running example: Angelucci et al. [2015], a randomized controlled trial study of the efficacy of microcredit based on 16,560 data points.

We can reverse the studies qualitative conclusions by removing 15 observations ($< 0.1\%$ of the data).

How do we find sets of influential points? Difficult in general!

We provide a **automatic approximation** with finite-sample guarantees.

Studying the approximation reveals the causes of non-robustness.

Dropping data: Mexico Microcredit

Consider Angelucci et al. [2015], a randomized controlled trial study of the efficacy of microcredit in Mexico based on 16,560 data points.

The variable “Beta” estimates the effect of microcredit in US dollars.

	Beta (SE)
Original result	-4.55 (5.88)

Original conclusion:

There is no evidence that microcredit is effective.

Dropping data: Mexico Microcredit

Consider Angelucci et al. [2015], a randomized controlled trial study of the efficacy of microcredit in Mexico based on 16,560 data points.

The variable “Beta” estimates the effect of microcredit in US dollars.

	Left out points	Beta (SE)
Original result	0	-4.55 (5.88)
Change sign	1	0.4 (3.19)

Original conclusion:

There is no evidence that microcredit is effective.

Dropping data: Mexico Microcredit

Consider Angelucci et al. [2015], a randomized controlled trial study of the efficacy of microcredit in Mexico based on 16,560 data points.

The variable “Beta” estimates the effect of microcredit in US dollars.

	Left out points	Beta (SE)
Original result	0	-4.55 (5.88)
Change sign	1	0.4 (3.19)
Change significance	14	-10.96 (5.57)

Original conclusion:

There is no evidence that microcredit is effective.

Dropping data: Mexico Microcredit

Consider Angelucci et al. [2015], a randomized controlled trial study of the efficacy of microcredit in Mexico based on 16,560 data points.

The variable “Beta” estimates the effect of microcredit in US dollars.

	Left out points	Beta (SE)
Original result	0	-4.55 (5.88)
Change sign	1	0.4 (3.19)
Change significance	14	-10.96 (5.57)
Change sign and significance	15	7.03 (2.55)

Original conclusion:

There is no evidence that microcredit is effective.

Dropping data: Mexico Microcredit

Consider Angelucci et al. [2015], a randomized controlled trial study of the efficacy of microcredit in Mexico based on 16,560 data points.

The variable “Beta” estimates the effect of microcredit in US dollars.

	Left out points	Beta (SE)
Original result	0	-4.55 (5.88)
Change sign	1	0.4 (3.19)
Change significance	14	-10.96 (5.57)
Change sign and significance	15	7.03 (2.55)

Original conclusion:

There is no evidence that microcredit is effective.

Potential conclusions after data dropping:

The effect of microcredit is positive (negative) & statistically significant.

Dropping data: Mexico Microcredit

Consider Angelucci et al. [2015], a randomized controlled trial study of the efficacy of microcredit in Mexico based on 16,560 data points.

The variable “Beta” estimates the effect of microcredit in US dollars.

	Left out points	Beta (SE)
Original result	0	-4.55 (5.88)
Change sign	1	0.4 (3.19)
Change significance	14	-10.96 (5.57)
Change sign and significance	15	7.03 (2.55)

Original conclusion:

There is no evidence that microcredit is effective.

Potential conclusions after data dropping:

The effect of microcredit is positive (negative) & statistically significant.

The culprit is signal to noise ratio.

By the end of the talk, we will see that the sensitivity is due to

- High variability of the outcome (household profit) relative to
- A small signal driving the conclusion (statistical significance)

Dropping data: Motivation

Would you be concerned if you could **reverse your conclusion** by removing a **small proportion** (say, 0.1%) of your data?

Dropping data: Motivation

Would you be concerned if you could **reverse your conclusion** by removing a **small proportion** (say, 0.1%) of your data?

Not always! But sometimes, surely yes.

Thinking without random noise can be helpful.

Suppose you have a farm, and want to know whether your average yield is greater than 170 bushels per acre. At harvest, you measure 200 bushels per acre.

- Scenario one: If your yield is greater than 170 bushels per acre, you make a profit.
 - Don't care about sensitivity to small subsets
- Scenario two: You want to recommend your farming methods to a friend across the valley.
 - Might care about sensitivity to small subsets

For example, often in economics:

- Small fractions of data are missing not-at-random,
- Policy population is different from analyzed population,
- We report a convenient summary (e.g. mean) of a complex effect,
- Models are stylized proxies of reality.

Which estimators do we study?

Z-estimators. Suppose we have N data points $\vec{d} = d_1, \dots, d_N$. Then:

$$\hat{\theta} := \vec{\theta} \text{ such that } \sum_{n=1}^N G(\vec{\theta}, d_n) = 0_P.$$

Examples: MLE, OLS, VB, &c (all minimizers of smooth empirical loss).

Function of interest. Qualitative decision based on $\phi(\hat{\theta}) \in \mathbb{R}$. E.g.:

- A particular component: $\phi(\theta) = \theta_d$
- The end of a confidence interval: $\phi(\theta) = \theta_d + \frac{1.96}{\sqrt{N}} \hat{\sigma}(\hat{\theta})$

Fix a proportion $0 < \alpha \ll 1$ of points to drop and find a set $\mathcal{S} \subset \{1, \dots, N\}$ with $|\mathcal{S}| \leq \lfloor \alpha N \rfloor$ that extremizes $\phi(\hat{\theta})$ when dropped.

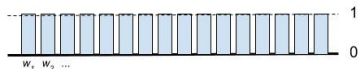
- **Problem:** There are many sets with $|\mathcal{S}| \leq \lfloor \alpha N \rfloor$.
 - E.g., in Angelucci et al. [2015], $\binom{16,560}{15} \approx 1.5 \cdot 10^{51}$
- **Problem:** Evaluating $\phi(\hat{\theta}(\vec{d}_{-\mathcal{S}}))$ requires an estimation problem.
 - E.g., in Angelucci et al. [2015] computing the OLS estimator.
 - Other examples are even harder (VB, machine learning)

An approximation is needed!

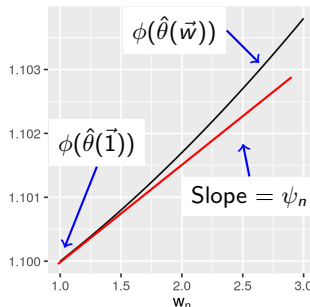
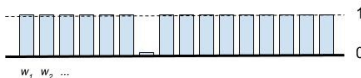
Which estimators do we study?

$$\hat{\theta} := \vec{\theta} \text{ such that } \sum_{n=1}^N G(\vec{\theta}, d_n) = 0_P.$$

Original weights: $\vec{1} = (1, \dots, 1)$



Leave points out by setting their elements of \vec{w} to zero.



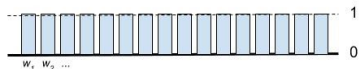
The slopes $\psi_n := \left. \frac{\partial \phi(\hat{\theta}(\vec{w}))}{\partial w_n} \right|_{\vec{1}}$ are values of the **empirical influence function** [Hampel, 1986]. We call them “influence scores.”

Second-order derivatives control the error of the linear approximation.

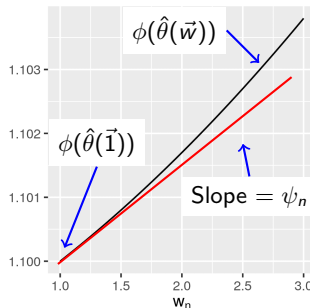
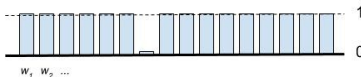
Which estimators do we study?

$$\hat{\theta}(\vec{w}) := \vec{\theta} \text{ such that } \sum_{n=1}^N \vec{w}_n G(\vec{\theta}, d_n) = 0_P.$$

Original weights: $\vec{1} = (1, \dots, 1)$



Leave points out by setting their elements of \vec{w} to zero.



The slopes $\psi_n := \left. \frac{\partial \phi(\hat{\theta}(\vec{w}))}{\partial w_n} \right|_{\vec{1}}$ are values of the **empirical influence function** [Hampel, 1986]. We call them “influence scores.”

Second-order derivatives control the error of the linear approximation.

Taylor series approximation.

Problem: How large can you make $\phi(\hat{\theta}(\vec{w}))$ leaving out no more than $\lfloor \alpha N \rfloor$ points? **Combinatorially hard!**

To simplify the search over \vec{w} , we form the Taylor series approximation:

$$\phi(\hat{\theta}(\vec{w})) \approx \phi^{\text{lin}}(\vec{w}) := \phi(\hat{\theta}(\vec{1})) + \sum_{n=1}^N \psi_n(\vec{w}_n - 1)$$

Approximate solution: How large can you make $\phi^{\text{lin}}(\vec{w})$ leaving out no more than $\lfloor \alpha N \rfloor$ points? **Easy!**

The most influential points for $\phi^{\text{lin}}(\vec{w})$ have the most negative ψ_n .

We provide **finite-sample theory** showing that

$$\left| \phi(\hat{\theta}(\vec{w})) - \phi^{\text{lin}}(\vec{w}) \right| = O \left(\left\| \frac{1}{N}(\vec{w} - \vec{1}) \right\|_2^2 \right) = O(\alpha) \text{ as } \alpha \rightarrow 0.$$

How to compute the influence scores ψ_n ?

By the chain rule, $\psi_n = \left. \frac{\partial \phi(\hat{\theta}(\vec{w}))}{\partial \vec{w}_n} \right|_{\vec{1}} = \left. \frac{d\phi(\theta)}{d\theta^T} \right|_{\hat{\theta}} \left. \frac{\partial \hat{\theta}(\vec{w})}{\partial \vec{w}_n} \right|_{\vec{1}}.$

Recall that $\hat{\theta}(\vec{w}) := \vec{\theta}$ such that $\sum_{n=1}^N \vec{w}_n G(\vec{\theta}, d_n) = 0_P.$

The **implicit function theorem** expresses $\left. \frac{\partial \hat{\theta}(\vec{w})}{\partial \vec{w}_n} \right|_{\vec{1}}$ as a linear system.

Computation of ψ_n is fully automatable from a software implementation of $G(\cdot, \cdot)$ and $\phi(\cdot)$ with **automatic differentiation** [Baydin et al., 2017].

We have an R package, `rgiordan/zaminfluence`, for OLS and IV.

Taylor series approximation.

Procedure:

Taylor series approximation.

Procedure:

- 1 Compute the “original” estimator, $\hat{\theta}(\vec{1})$ and $\phi(\hat{\theta}(\vec{1}))$.

Procedure:

- 1 Compute the “original” estimator, $\hat{\theta}(\vec{1})$ and $\phi(\hat{\theta}(\vec{1}))$.
- 2 Let Δ denote an increase in $\phi(\hat{\theta})$ that would change conclusions.

Procedure:

- 1 Compute the “original” estimator, $\hat{\theta}(\vec{1})$ and $\phi(\hat{\theta}(\vec{1}))$.
- 2 Let Δ denote an increase in $\phi(\hat{\theta})$ that would change conclusions.
- 3 Compute and sort the influence scores, $\psi_{(1)} \leq \psi_{(2)} \leq \dots \leq \psi_{(N)}$.

Procedure:

- 1 Compute the “original” estimator, $\hat{\theta}(\vec{1})$ and $\phi(\hat{\theta}(\vec{1}))$.
- 2 Let Δ denote an increase in $\phi(\hat{\theta})$ that would change conclusions.
- 3 Compute and sort the influence scores, $\psi_{(1)} \leq \psi_{(2)} \leq \dots \leq \psi_{(N)}$.
- 4 Let \vec{w}^* leave out the data corresponding to $\psi_{(1)}, \dots, \psi_{(\lfloor \alpha N \rfloor)}$.

Procedure:

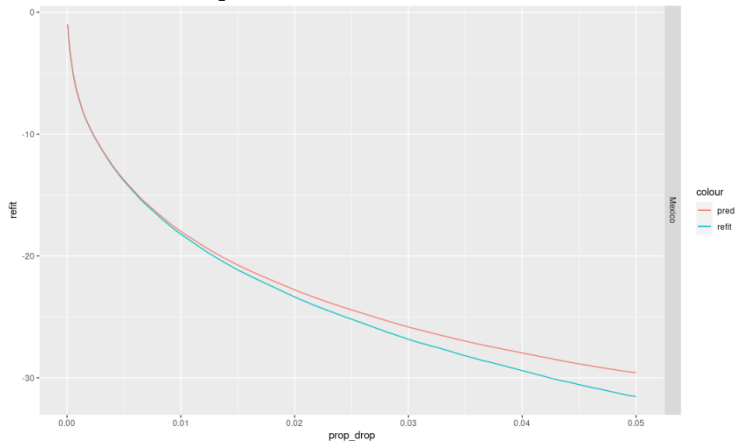
- 1 Compute the “original” estimator, $\hat{\theta}(\vec{1})$ and $\phi(\hat{\theta}(\vec{1}))$.
- 2 Let Δ denote an increase in $\phi(\hat{\theta})$ that would change conclusions.
- 3 Compute and sort the influence scores, $\psi_{(1)} \leq \psi_{(2)} \leq \dots \leq \psi_{(N)}$.
- 4 Let \vec{w}^* leave out the data corresponding to $\psi_{(1)}, \dots, \psi_{(\lfloor \alpha N \rfloor)}$.
- 5 Report non-robustness if $\phi^{\text{lin}}(\vec{w}^*) - \phi(\hat{\theta}) = -\sum_{n=1}^{\lfloor \alpha N \rfloor} \psi_{(n)} \geq \Delta$.

Procedure:

- 1 Compute the “original” estimator, $\hat{\theta}(\vec{1})$ and $\phi(\hat{\theta}(\vec{1}))$.
- 2 Let Δ denote an increase in $\phi(\hat{\theta})$ that would change conclusions.
- 3 Compute and sort the influence scores, $\psi_{(1)} \leq \psi_{(2)} \leq \dots \leq \psi_{(N)}$.
- 4 Let \vec{w}^* leave out the data corresponding to $\psi_{(1)}, \dots, \psi_{(\lfloor \alpha N \rfloor)}$.
- 5 Report non-robustness if $\phi^{\text{lin}}(\vec{w}^*) - \phi(\hat{\theta}) = -\sum_{n=1}^{\lfloor \alpha N \rfloor} \psi_{(n)} \geq \Delta$.
- 6 **Optional:** Compute $\hat{\theta}(\vec{w}^*)$, and verify that $\phi(\hat{\theta}(\vec{w}^*)) - \phi(\hat{\theta}) \geq \Delta$.

Mexico example:

See `microcredit_profit_sandbox.R`.



Selected experimental results.

- The “Refit estimate” column shows the result of re-fitting the model removing the points with the largest influence scores.
- Stars indicate significance at the 5% level.
- Refits that achieved the desired change are bolded.

Study case	Original estimate	Target change	Refit estimate	Observations dropped
Mexico	-4.549 (5.879)	Sign change	0.398 (3.194)	1 = 0.01%
		Significance change	-10.962 (5.565)*	14 = 0.08%
		Significant sign change	7.030 (2.549)*	15 = 0.09%

Table: Microcredit Mexico results [Angelucci et al., 2015].

Selected experimental results.

- The “Refit estimate” column shows the result of re-fitting the model removing the points with the largest influence scores.
- Stars indicate significance at the 5% level.
- Refits that achieved the desired change are bolded.

Study case	Original estimate	Target change	Refit estimate	Observations dropped
Mexico	-4.549 (5.879)	Sign change	0.398 (3.194)	1 = 0.01%
		Significance change	-10.962 (5.565)*	14 = 0.08%
		Significant sign change	7.030 (2.549)*	15 = 0.09%

Table: Microcredit Mexico results [Angelucci et al., 2015].

Study case	Original estimate	Target change	Refit estimate	Observations dropped
Poor, period 10	33.861 (4.468)*	Sign change	-2.559 (3.541)	697 = 6.63%
		Significance change	4.806 (3.684)	435 = 4.14%
		Significant sign change	-9.416 (3.296)*	986 = 9.37%
Non-poor, period 10	21.493 (9.405)*	Sign change	-0.573 (6.750)	30 = 0.70%
		Significance change	16.262 (8.927)	3 = 0.07%
		Significant sign change	-10.845 (6.467)	92 = 2.16%

Table: Cash transfers results. [Angelucci and De Giorgi, 2009]

A simulation

For $N = 5,000$ data points, compute the OLS estimator from:

Regressors
 $x_n \sim \mathcal{N}(0, \sigma_x^2)$

Residuals
 $\varepsilon_n \sim \mathcal{N}(0, \sigma_\varepsilon^2)$

Responses
 $y_n = 0.5x_n + \varepsilon_n$

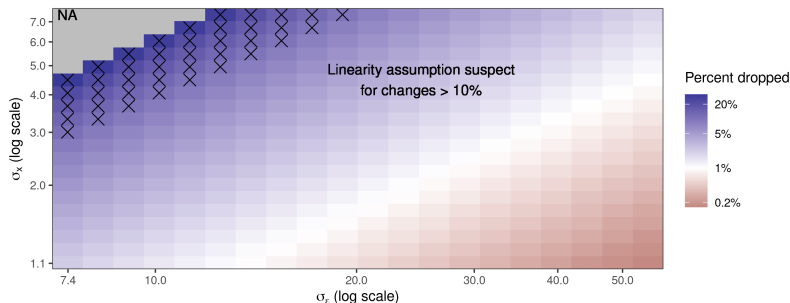


Figure: The approximate perturbation inducing proportion at differing values of σ_x and σ_ε . Red colors indicate datasets whose sign can be predicted to change when dropping less than 1% of datapoints. The grey areas indicate $\hat{\Psi}_\alpha = \text{NA}$, a failure of the linear approximation to locate any way to change the sign.

What makes an estimator non-robust? A tail sum.

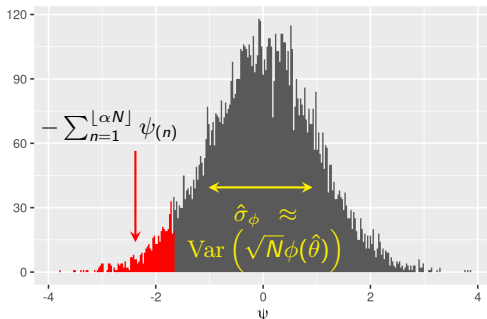
Report non-robustness if:

$$\Delta \leq \phi^{\text{lin}}(\vec{w}^*) - \phi(\hat{\theta}) = - \sum_{n=1}^{\lfloor \alpha N \rfloor} \psi_{(n)} =: \hat{\sigma}_{\phi} \hat{\mathcal{J}}_{\alpha}$$

We will show that:

- The “noise” $\hat{\sigma}_{\phi}^2 \rightarrow \text{Var}(\sqrt{N}\phi)$ [Hampel, 1986]
- The “shape” $\hat{\mathcal{J}}_{\alpha} \leq \sqrt{\alpha(1-\alpha)}$ and converges to a nonzero constant

Influence score histogram (N = 10000, $\alpha = 0.05$)



Corollaries.

Report non-robustness if:

$$\Delta \leq \phi^{\text{lin}}(\vec{w}^*) - \phi(\hat{\theta}) = \hat{\sigma}_\phi \hat{\mathcal{J}}_\alpha \quad \Leftrightarrow \quad \frac{\Delta}{\hat{\sigma}_\phi} \leq \hat{\mathcal{J}}_\alpha.$$

We call $\frac{\Delta}{\hat{\sigma}_\phi}$ the “signal to noise ratio.”

Report non-robustness if:

$$\Delta \leq \phi^{\text{lin}}(\vec{w}^*) - \phi(\hat{\theta}) = \hat{\sigma}_{\phi} \hat{\mathcal{J}}_{\alpha} \quad \Leftrightarrow \quad \frac{\Delta}{\hat{\sigma}_{\phi}} \leq \hat{\mathcal{J}}_{\alpha}.$$

We call $\frac{\Delta}{\hat{\sigma}_{\phi}}$ the “signal to noise ratio.”

Corollary: Non-robustness possible even with correct specification.

Report non-robustness if:

$$\Delta \leq \phi^{\text{lin}}(\vec{w}^*) - \phi(\hat{\theta}) = \hat{\sigma}_{\phi} \hat{\mathcal{J}}_{\alpha} \quad \Leftrightarrow \quad \frac{\Delta}{\hat{\sigma}_{\phi}} \leq \hat{\mathcal{J}}_{\alpha}.$$

We call $\frac{\Delta}{\hat{\sigma}_{\phi}}$ the “signal to noise ratio.”

Corollary: Non-robustness possible even with correct specification.

Corollary: Leave- $\lfloor \alpha N \rfloor$ -out robustness does not vanish as $N \rightarrow \infty$.

Report non-robustness if:

$$\Delta \leq \phi^{\text{lin}}(\vec{w}^*) - \phi(\hat{\theta}) = \hat{\sigma}_\phi \hat{\mathcal{J}}_\alpha \quad \Leftrightarrow \quad \frac{\Delta}{\hat{\sigma}_\phi} \leq \hat{\mathcal{J}}_\alpha.$$

We call $\frac{\Delta}{\hat{\sigma}_\phi}$ the “signal to noise ratio.”

Corollary: Non-robustness possible even with correct specification.

Corollary: Leave- $\lfloor \alpha N \rfloor$ -out robustness does not vanish as $N \rightarrow \infty$.

Recall that standard errors reject when $\frac{\Delta}{\hat{\sigma}_\phi} \leq \frac{1.96}{\sqrt{N}}$.

Report non-robustness if:

$$\Delta \leq \phi^{\text{lin}}(\vec{w}^*) - \phi(\hat{\theta}) = \hat{\sigma}_{\phi} \hat{\mathcal{J}}_{\alpha} \quad \Leftrightarrow \quad \frac{\Delta}{\hat{\sigma}_{\phi}} \leq \hat{\mathcal{J}}_{\alpha}.$$

We call $\frac{\Delta}{\hat{\sigma}_{\phi}}$ the “signal to noise ratio.”

Corollary: Non-robustness possible even with correct specification.

Corollary: Leave- $\lfloor \alpha N \rfloor$ -out robustness does not vanish as $N \rightarrow \infty$.

Recall that standard errors reject when $\frac{\Delta}{\hat{\sigma}_{\phi}} \leq \frac{1.96}{\sqrt{N}}$.

Corollary: Leave- $\lfloor \alpha N \rfloor$ -out is different from standard errors.

Report non-robustness if:

$$\Delta \leq \phi^{\text{lin}}(\vec{w}^*) - \phi(\hat{\theta}) = \hat{\sigma}_\phi \hat{\mathcal{J}}_\alpha \quad \Leftrightarrow \quad \frac{\Delta}{\hat{\sigma}_\phi} \leq \hat{\mathcal{J}}_\alpha.$$

We call $\frac{\Delta}{\hat{\sigma}_\phi}$ the “signal to noise ratio.”

Corollary: Non-robustness possible even with correct specification.

Corollary: Leave- $\lfloor \alpha N \rfloor$ -out robustness does not vanish as $N \rightarrow \infty$.

Recall that standard errors reject when $\frac{\Delta}{\hat{\sigma}_\phi} \leq \frac{1.96}{\sqrt{N}}$.

Corollary: Leave- $\lfloor \alpha N \rfloor$ -out is different from standard errors.

Corollary: Insignificance is always non-robust.

Take $\Delta = \frac{1.96 \hat{\sigma}_\phi}{\sqrt{N}} \rightarrow 0 \leq \hat{\mathcal{J}}_\alpha$.

Report non-robustness if:

$$\Delta \leq \phi^{\text{lin}}(\vec{w}^*) - \phi(\hat{\theta}) = \hat{\sigma}_{\phi} \hat{\mathcal{J}}_{\alpha} \quad \Leftrightarrow \quad \frac{\Delta}{\hat{\sigma}_{\phi}} \leq \hat{\mathcal{J}}_{\alpha}.$$

We call $\frac{\Delta}{\hat{\sigma}_{\phi}}$ the “signal to noise ratio.”

Corollary: Non-robustness possible even with correct specification.

Corollary: Leave- $\lfloor \alpha N \rfloor$ -out robustness does not vanish as $N \rightarrow \infty$.

Recall that standard errors reject when $\frac{\Delta}{\hat{\sigma}_{\phi}} \leq \frac{1.96}{\sqrt{N}}$.

Corollary: Leave- $\lfloor \alpha N \rfloor$ -out is different from standard errors.

Corollary: Insignificance is always non-robust.

Take $\Delta = \frac{1.96 \hat{\sigma}_{\phi}}{\sqrt{N}} \rightarrow 0 \leq \hat{\mathcal{J}}_{\alpha}$.

Corollary: Gross outliers primarily affect robustness through $\hat{\sigma}_{\phi}$.

Cauchy-Schwartz is tight when all the influence scores are the same.

Conclusion:

Related work and future directions

Tamara Broderick, Ryan Giordano, Rachael Meager (alphabetical authors)
“An Automatic Finite-Sample Robustness Metric: Can Dropping a Little Data Change Conclusions?”

<https://arxiv.org/abs/2011.14999>

-
- M. Angelucci and G. De Giorgi. Indirect effects of an aid program: How do cash transfers affect ineligibles' consumption? *American Economic Review*, 99(1):486–508, 2009.
- M. Angelucci, D. Karlan, and J. Zinman. Microcredit impacts: Evidence from a randomized microcredit program placement experiment by Compartamos Banco. *American Economic Journal: Applied Economics*, 7(1):151–82, 2015.
- A. Baydin, B. Pearlmutter, A. Radul, and J. Siskind. Automatic differentiation in machine learning: A survey. *The Journal of Machine Learning Research*, 18(1):5595–5637, 2017.
- F. Hampel. *Robust statistics: The approach based on influence functions*, volume 196. Wiley-Interscience, 1986.